

Citation for published version:

Dorta, G, Vicente, S, Agapito, L, Campbell, N & Simpson, I 2018, 'Structured Uncertainty Prediction Networks'
Paper presented at IEEE Conference on Computer Vision and Pattern Recognition (CVPR) 2018, 18/06/18 -
22/06/18, .

Publication date:

2018

Document Version

Publisher's PDF, also known as Version of record

[Link to publication](#)

University of Bath

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Supplement: Structured Uncertainty Prediction Networks

Garoe Dorta^{1,2} Sara Vicente² Lourdes Agapito³ Neill D.F. Campbell¹ Ivor Simpson²

¹University of Bath ²Anthropics Technology Ltd. ³University College London

¹{g.dorta.perez,n.campbell}@bath.ac.uk ²{sara,ivor}@anthropics.com ³l.agapito@cs.ucl.ac.uk

A. Network architectures

All the models were trained with a batch size of 64. The exp block in all the architectures removes the log in the diagonal values of the Cholesky matrix that is being estimated.

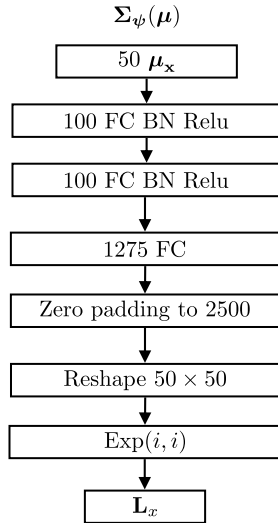


Figure 1: Covariance estimation network architecture for the splines dataset.

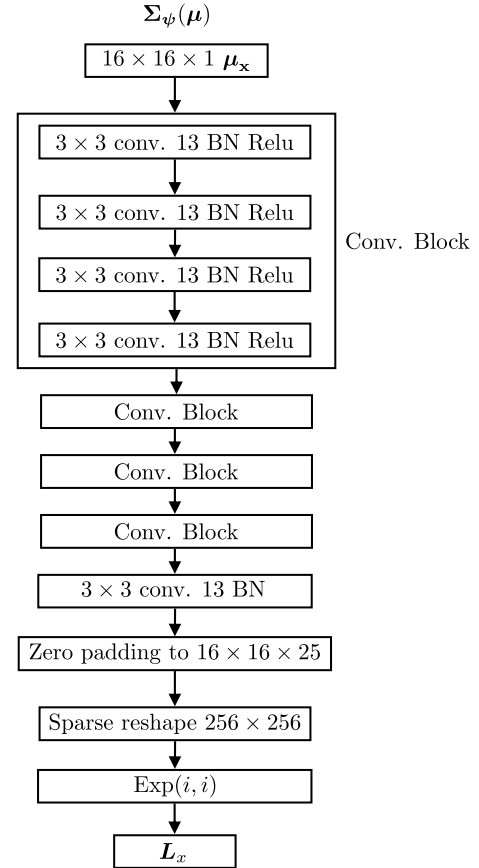


Figure 2: Covariance estimation network architecture for the ellipses dataset.

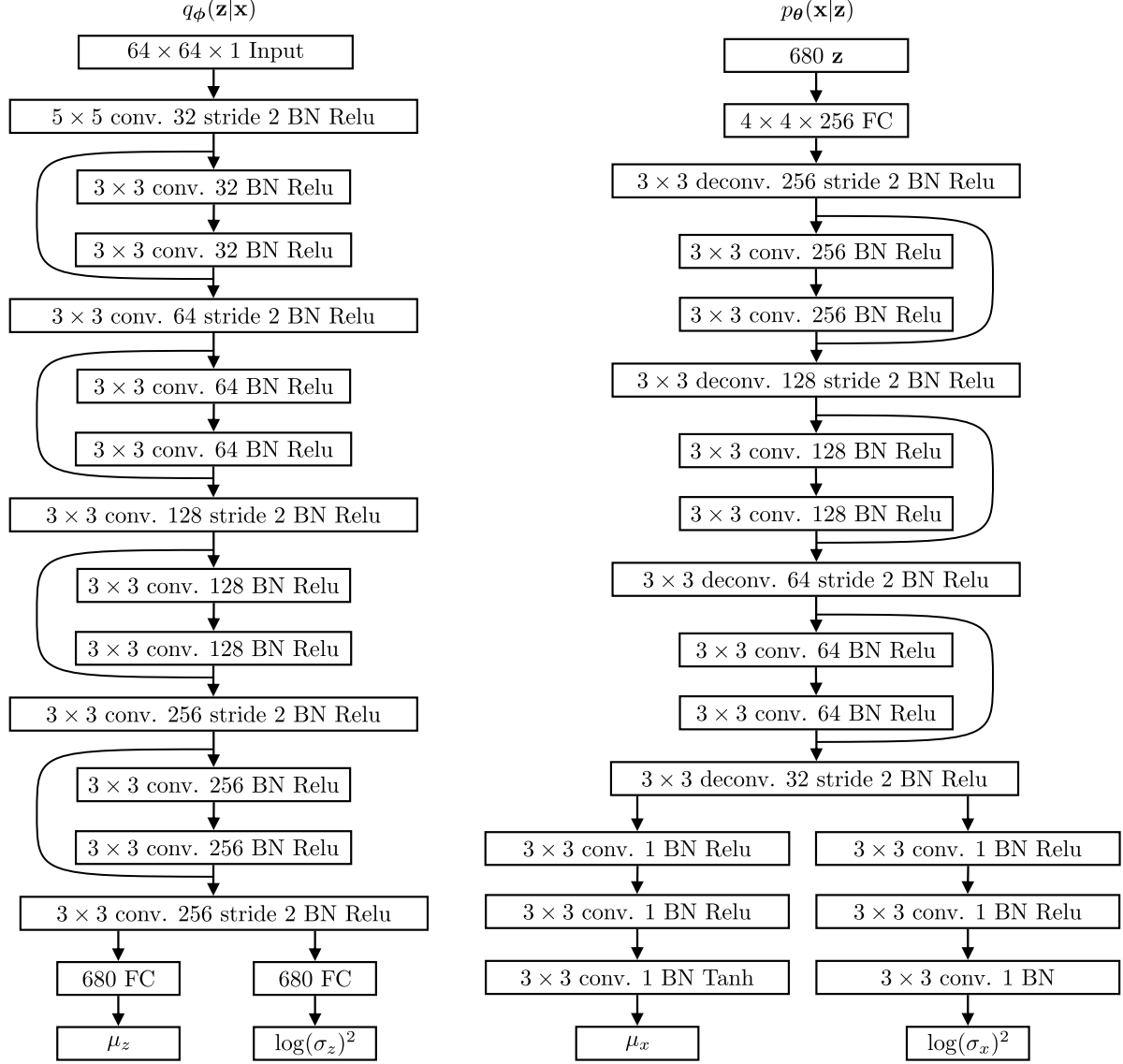


Figure 3: VAE architecture for the CelebA dataset.

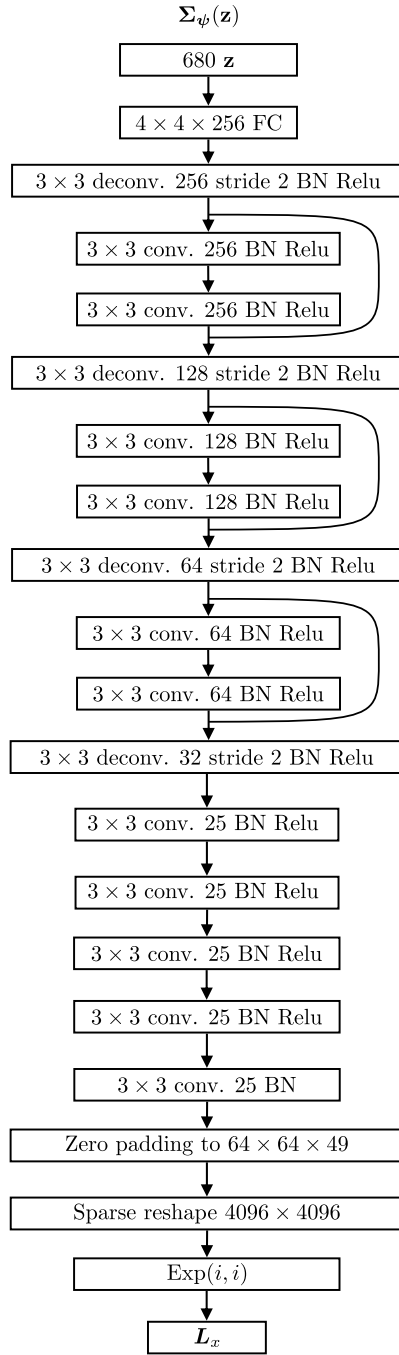


Figure 4: Covariance estimation network architecture for the CelebA dataset.

B. CIFAR10 dataset

In this dataset our models uses a patch size of 5×5 , and an architecture adapted from the CelebA model for 32×32 images. Quantitatively, a VAE with diagonal covariance achieved a marginal negative log likelihood of -1026 ± 462 and our model of -1333 ± 477 , where the likelihood was evaluated with 500 z samples per image.

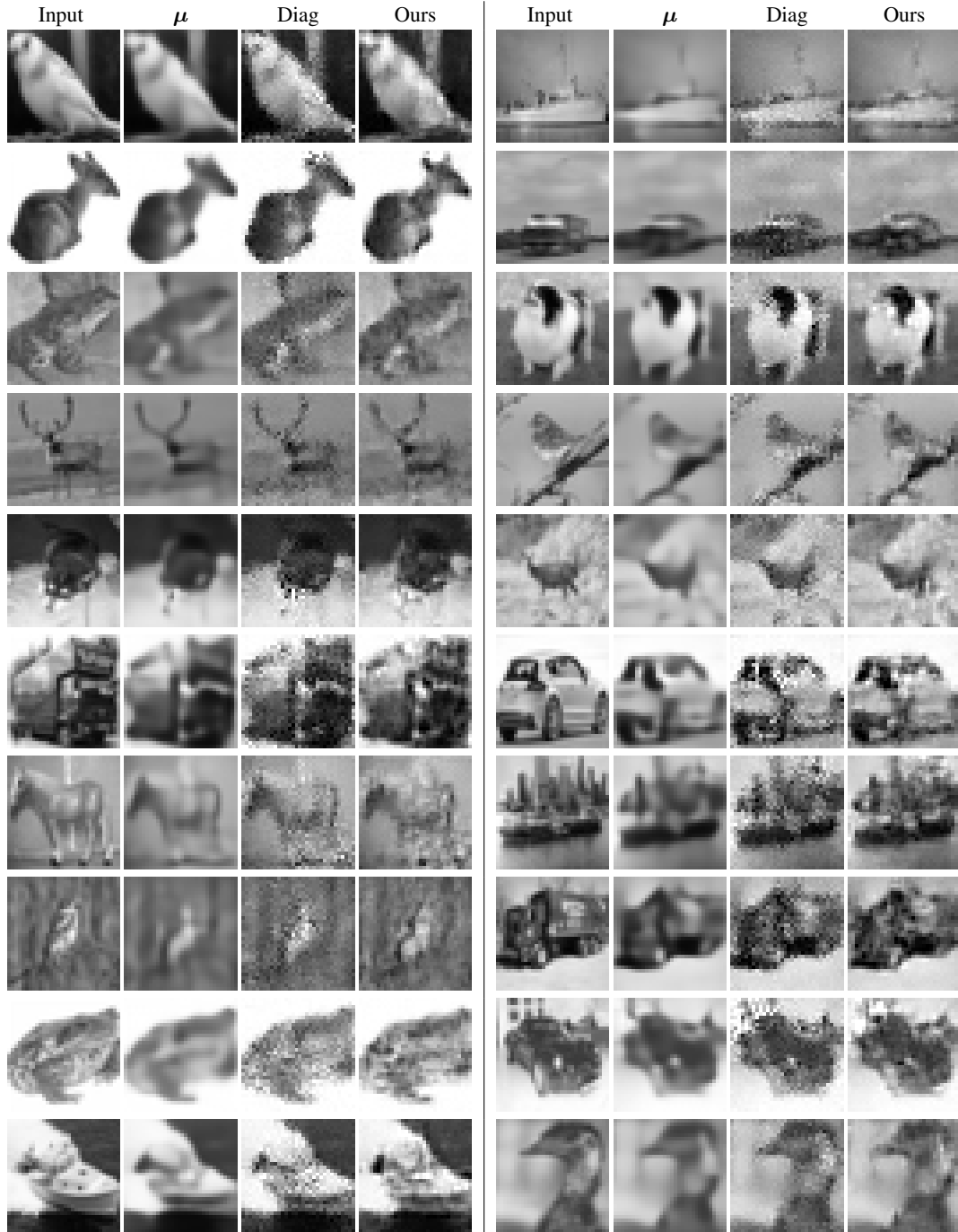


Figure 5: Comparison of image reconstructions on the CIFAR10 dataset, where μ is a reconstruction from a VAE, Diag. is $\mu + \epsilon$, where ϵ is a sample from a diagonal covariance matrix, and in our method ϵ is a sample from a dense covariance matrix.

C. Splines dataset

In the splines dataset, for each example, μ is a cubic spline interpolation of five random points sampled from a Gaussian distribution, as shown in Fig. 6c. Given a predefined prototype covariance matrix (shown in Fig. 6a), the covariance matrix Σ for a particular example is constructed by scaling this prototype covariance by the absolute value of μ , as shown in Fig. 6b. Finally, a random sample is drawn from the covariance matrix and added to the mean μ , which generates the final example $\mathbf{x} = \mu + \epsilon$, as shown in Fig. 6d.

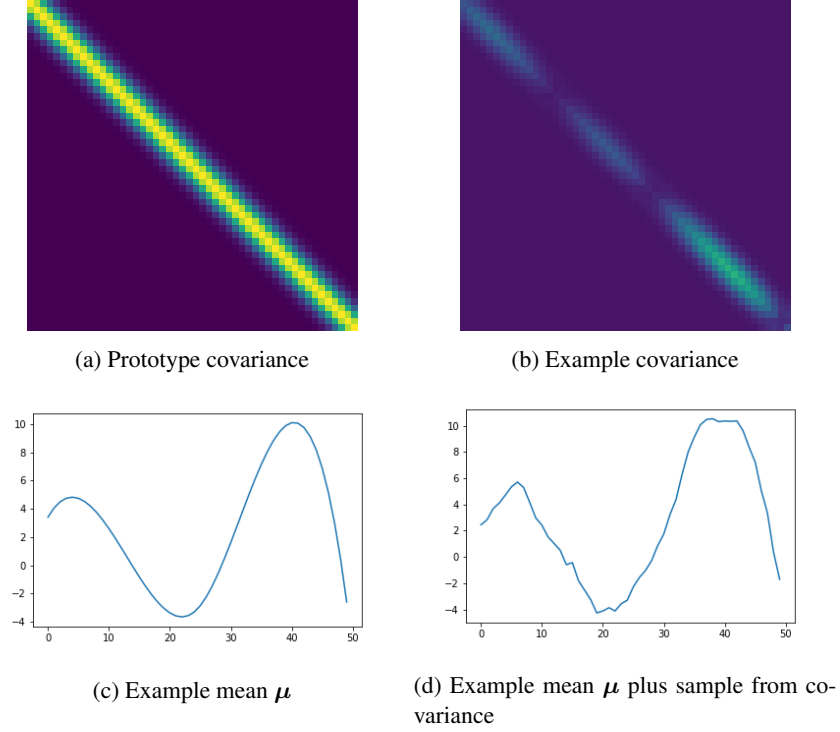


Figure 6: The splines dataset is synthesized using a prototype covariance matrix (a), that is transformed to generate a number of example covariances (b), where each transformation is function of a different spline (c). Each example (d) in the dataset is constructed by taking a single sample from a covariance (b) and adding it to the corresponding spline (c).

D. CelebA



Figure 7: Denoising experiment, left column: original image without noise, second column: image with added noise, third column: denoising autoencoder (DAE) result, fourth column: our result, fifth column: VAE reconstruction from the noisy input, sixth column: difference between the VAE reconstruction and the noisy input, right column: the difference projected on $\hat{\Sigma}$, the matrix constructed with 1000 eigenvectors of Σ with the largest eigenvalues. Our result is the sum of the projected difference and the VAE reconstruction. Our model is able to recover fine details that are lost with the DAE approach.





Figure 9: Comparison of image reconstructions for the different models. The best of 100 and 1000 samples from Σ measured using the MSE to the input are shown. The AE and VAE both generate over-smoothed images. For both the AE and VAE, our model adds plausible high-frequencies from a single sample drawn from the predicted uncertainty distribution.



Figure 10: Images generated by decoding samples from the prior distribution on the latent space of a β -VAE with added residuals from our model, where $\beta = 5$.



Figure 11: Images generated by decoding samples from the prior distribution on the latent space of a VAE.

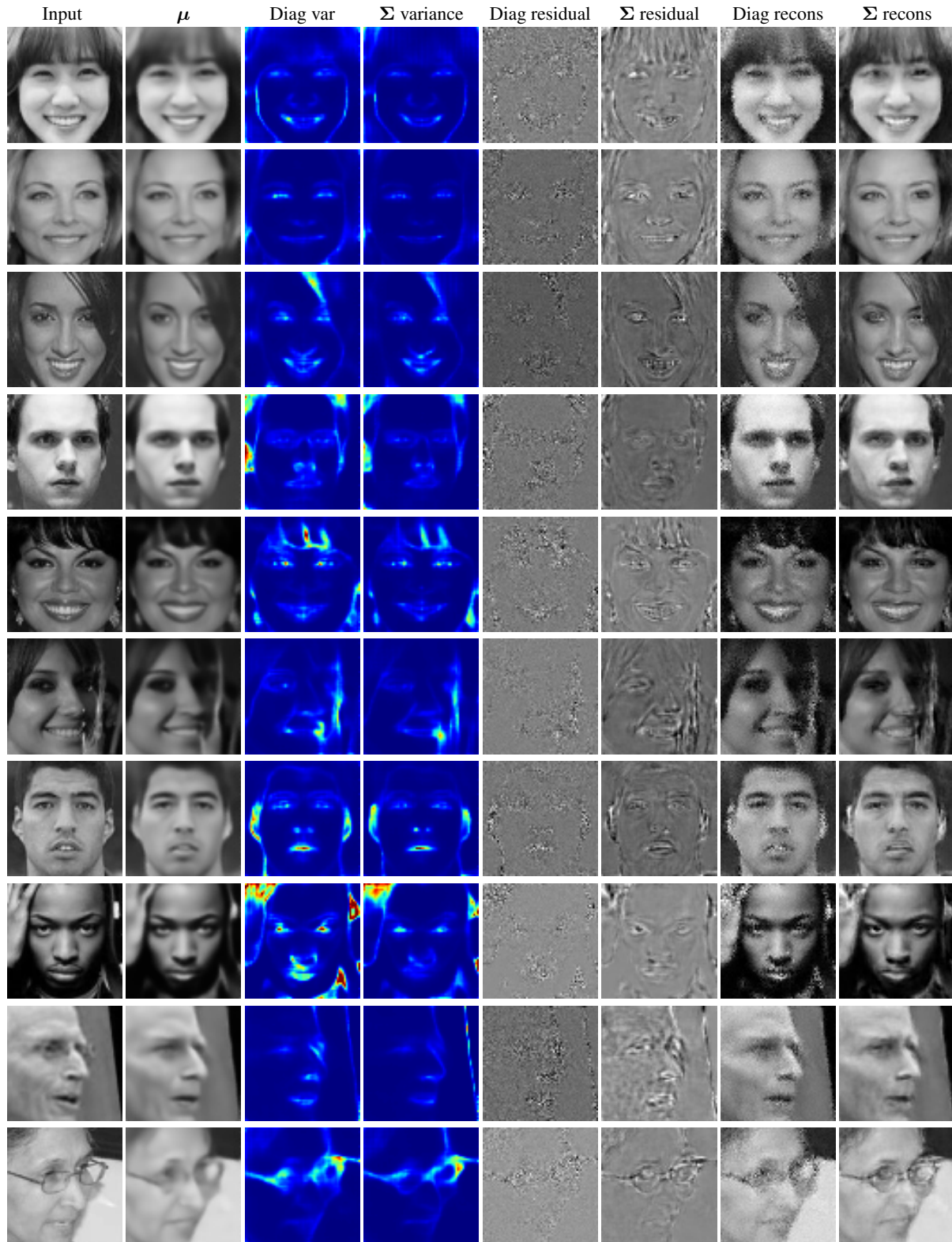


Figure 12: Variance maps for different inputs, taking a single sample from the estimated residual distribution. The diagonal noise estimation model mistakenly identifies teeth or skin wrinkles as variance, whereas the covariance model properly identifies them as regions with high covariance, yet low variance.

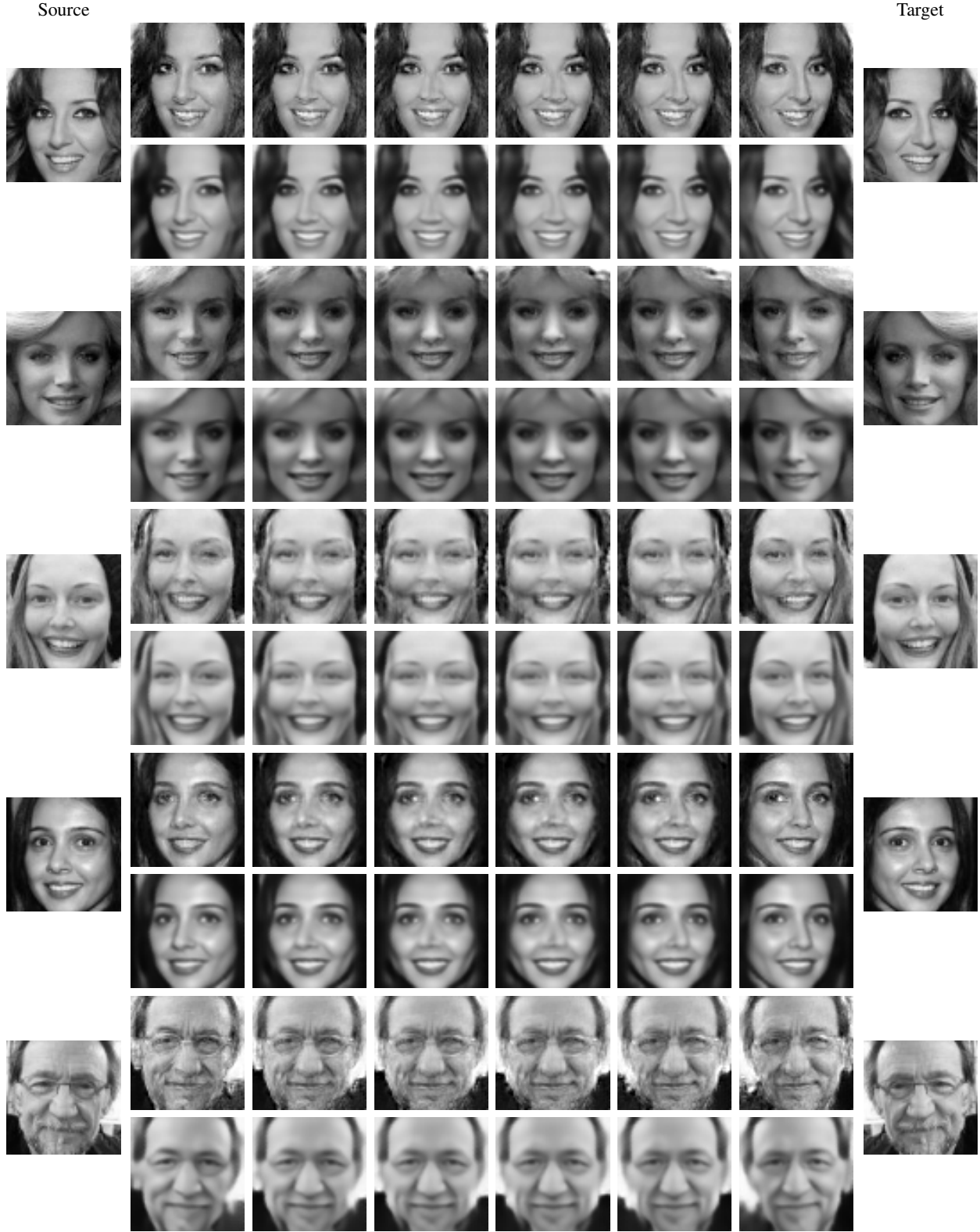


Figure 13: Samples drawn with our model while interpolating on the latent space, from the source to target. Using a fixed noise vector \mathbf{u} , a single sample is drawn from our model as $\mathbf{x} = \boldsymbol{\mu} + \mathbf{M}\mathbf{u}$, where the covariance $\boldsymbol{\Sigma} = \mathbf{M}\mathbf{M}^T$ is estimated from each \mathbf{z} . The latent values \mathbf{z} are interpolated using polar interpolation. For every pair of rows, the first contains the result from our model, and the second contains the reconstruction $\boldsymbol{\mu}(\mathbf{z})$ from the VAE.